



A SARS vírus genetikai állományának megfejtése

2003. április 7-én 1 órakor érkezett meg a SARS vírus tenyészete a Michael Smith Genomtudományi Központba. Öt nappal később a laboratórium elsőként hozta nyilvánosságra a vírus genetikai állományának nukleotidsorrendjét.

dén áprilisban a Genomtudományi Központban (Genome Sciences Centre, azaz GSC) tettük közzé az első teljes gén-készletszerkezetét annak a koronavírusrnak, amely ismereteink szerint a Severe Acute Respiratory Syndrome (SARS) járvány okozója. A GSC-nél az 1999-es kezdetek óta minden vizsgálatot, tárolást és hálózati háttér Linux-rendszerek alatt végeznek. A SARS vírus projektben az adatok tárolását, feldolgozását és nyilvánosságra hozását számos Linux-kiszolgáló végezte, kezdve a pehelysúlyú, de hasznos IBM x330-tól egészen a behemót nyolcutas Xeon x440-ig. A Linux által nyújtott rugalmas háttér lehetővé tette, hogy a megfejtési folyamat szinte minden lépését automatizáljuk. A Linux-közösség támogatásával és a hírcsoportok, webes cikkek és HOGYAN-ok segítségével hihetetlenül olcsón munkára tudtuk fogni a középkeletű alkatrészeket.

A SARS első dokumentált megjelenése óta (2002. november 16.) a vírust összesen 8458 esetben észlelték Kínában (92%), Kanadában (3%), Szingapúrban (2%) és az Észak-Amerikai Egyesült Államokban (1%), valamint több mint 25 egyéb országban. A SARS halálozási esélye közelítőleg 5–10%, a 60 felettiiek esetében azonban 50% körüli. 2003. június 24-re a SARS már 807 életet követelt, igen mély negatív hatást gyakorolva az érintett régiók gazdaságára – egyedül Kína több milliárd dollárt veszített turisztikai és adójövedelmeiből.

2003. március 27-én *Marco Marra*, központunk igazgatója és *Caroline Astell*, projektünk vezetője úgy döntött, hogy megfejti a SARS koronavírus genomjának szekvenciáját. 2003. április 7-én éjjel 1 órakor egy torontói páciensből származó vírus, a Tor2 izolátum genetikai anyagának közel 50 ng-ja érkezett a kanadai Winnipeg 4. szintű Nemzeti Mikrobiológiai Laboratóriumából. Öt nappal később, 2003. április 12-én a Tor2 (Tor2/SARS) koronavírus genetikai állományának 29751 nukleotidhosszúságú

része már felkerült Apache kiszolgálónk Zope/Plone alapú oldalára – elérhetővé téve azt a teljes nyilvánosság számára. Néhány nappal később az úgynevezett Urbani izolátum szekvenciáját küldte el a CDC (Centers for Disease Control) Atlantából, Georgia államból.

A biológia virágzásnak indul

Az 1990-es évek előtt nem létezett olyan módszer, amellyel nagy mennyiségű nukleotidsorrend-adatot gyorsan meg lehetett volna határozni. Az Emberi Genom Projekt (Human Genome Project, azaz HGP) 1991-ben kezdődött, és 1999-re a sorrendnek mindössze 15 százalékát sikerült megfejteni. Ugyanakkor az 1990-es években kifejlesztett új módszereknek hála a HGP gyorsan közeledett a befejezés felé. 2000 közepe táján az emberi szekvencia kilencven százaléka már elérhető volt, és mostanra az emberi génállomány nukleotidsorrendje lényegében rendelkezésünkre áll. A HGP-hez hasonló projektek eredményei nyilvánosan is elérhetők az NCBI GenBank oldalain.

Működésének első tíz éve során (1982–1992) a Génbank valamivel több mint 100 MB-nyi szekvenciát gyűjtött össze 80 ezer bejegyzésben. A következő évtizedben (1992–2002) a Génbank rakétasebességgel növekedésnek indult, és az adatbázis elérte a 29 GB-ot – az emberi genom tízszeresét – 22 millió bejegyzésbe szedve. A Génbank minden nap tízezer bázissorrendadatot kap a világ különféle laborjaitól. Az egyik ilyen labor a GSC, amely 2003. április 13-án jelentette be a GenBankban a Tor2/SARS szekvenciáját. Ha kíváncsiak vagyunk, hogy milyen szerepet játszott a Linux abban a folyamatban, amely végül a GI:29826276 számú bejegyzés megszületéséhez vezetett, vissza kell nyúlunk a kezdetekig.

Parancssoros bioinformatika

Valószínűleg meg egy kis bioinformatikát a bash és néhány másik, a `/bin` és `/usr/bin` könyvtárban bujkáló program segítségével.

A Tor2/SARS genom GC arányát fogjuk kiszámítani – azaz a G-C vagy C-G bázispárok részarányát. Hogy érdekes legyen a dolog, az awk programot nem fogjuk használni. Először is wget-tel töltsük le a sorozatot, a `-q` kapcsolóval csillapítva a kimenetét:

```
> wget -q
↳ http://mkweb.bcgsc.ca/sars/AY274119.fa
> head AY274119.fa
gi|30248028|gb|AY274119.3| SARS coronavirus
↳ TOR2
ATATTTAGGTTTTTACCTACCCAGGA...
```

A nukleotidsorrend-fájlok FASTA formátumban vannak, amely a fejlécsort és magát a rögzített hosszúságú sorokra osztott nukleo-

tidláncot tartalmazza. A következő kód megszámlálja, hogy hány G és C található a láncban, majd az eredményt az összes bázis arányában jeleníti meg:

```
> grep -v ">" AY274119.fa | fold -w 1 |
tr "ATGC" "..xx" | sort | uniq -c |
sed 's/[^0-9]//g' | t -s "\012" " " |
sed 's/\([0-9]*\) \([0-9]*\) /scale = 3;
↳ \2 \ / (\1+\2) /' |
bc -i
scale = 3; 12127 / (17624+12127)
.407
```

Szekvenciánk 29 751 bázisból tehát 12 127 elem lesz akár G vagy C, így a GC-tartalom 41%-ra adódik.



1. kép Szekvenálólaborunk panorámája: 1. folyamatoknak megfelelő vonalkódok, 2. a Tango folyadékkezelő felület, 3. -80° C-os mélyhűtők, 4. áramforrások a PCR (polimeráz láncreakció) készülékekhez, 5. ABI 3730XL szekvenátorok, 6. ABI 3700 szekvenátorok, 7. x330 vezérlőfürt, 8. hálózati és áramcsatlakozók 9. a szekvenátorok ventilátorjára

0-18 TB három év alatt

1999 júniusában a labor hat szép bézsszínű számítógépet és közel ugyanannyi embert alkalmazott. A központi fájlkiszolgáló (2×Pentium 3, 400 MHz, 512 MB RAM, Red Hat 5.2 és 2.0.36-os rendszeremag) három RAID-0 18 GB SCSI-merevlemezelt kezelt DPT IV kártyán keresztül. Újabb 50 GB programozott RAID került a második gépbe (Pentium III, 400 MHz). További három Linux-ügyféllel és egy Microsoft Windows NT állomással együtt alkották a BC Cancer Agency (BCCA) hálózatát. Megszületésünk időpontja nagy előnyünkre szolgált. Mint minden kutatólabor, mi is lemezeket osztunk meg, folyamatokat kezelünk, programokat fordítunk, valamint adatokat tárolunk és kezelünk. Más szavakkal: éppen olyasmit csinálunk, amiben a Unix kiváló. Ha két-három évvel korábban kezdünk, az akkor még ifjú Linux bevezetése nem lett volna könnyű. Így ma valószínűleg ahelyett, hogy az olcsó kiöregedett PC-inket irodába számúzzuk vagy egyéb kevésbé nagyfokú hálózati feladatokra osztjuk be, megpróbálhatnánk a legjobb árat kapni az igen jelentős összegekbe került, már kiöregedett Sun kiszolgálóinkért. Szerencsére kiderült, hogy lehetőségünk van viszonylag olcsó PC-eket vásárolni, majd Linuxot telepítve rájuk nagyméretű, rugalmas és elképesztően költséghatékony Unix-környezetbe jutnunk. A Linuxnak köszönhetően többé már nem volt szükség rá, hogy egy ember fizetését Unix-munkaállomásokra költjük.

Éppen jó időben választottuk a Linuxot. A 2.0-s rendszeremag sziklaszilárd volt; az NFS kiszolgáló megerősödött, és teljes értékű asztali környezetek között válogathattunk. A létfontosságú bioinformatikai analízis-eszköz-készleteket letölthettük és lefordíthattuk. Ilyen például a nyílt forrású HGP: BLAST (sorrend-összehasonlító), a Phred (a szekvenátor által készített jelek értelmezése), a Phrap (sorrendek összeállítása) és a Consed (nukleotid- és aminosavsorrendek összeállításainak megjelenítése), továbbá néhány nukleinsav és fehérje-adatbázis. Természetesen a Perl volt a „mindenes” ezekben a műveletekben. A számítástechnikai munka elindításához igen kevés pénzt használtunk fel, így a nagy összegeket sokkal hatékonyabban költhettük a labor fejlesztésére (1. kép).

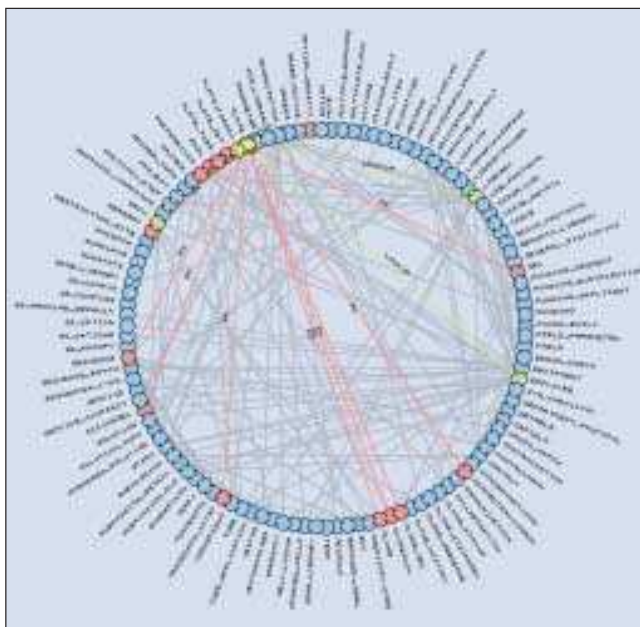
A Linux elcsípi a SARS-t

1999 őszén megkaptuk első automata DNS-szekvenátorunkat, egy MegaBACE 1000-est (6. kép). A szekvenátor segítségével egy DNS-mintában megállapítható a nukleotidok sorrendje, a módszer azonban jelenleg 5–800 bázis meghatározására korlátozódik egy időben. Ez az olvasási hossz jóval kisebb a jelenleg ismert legkisebb genomnál is (a Tor2/SARS mérete harmincezer nukleotid). Ezért az automata szekvenátor egyszerre 96 mintát kezel; vannak olyan típusok is, amelyekben egyszerre több, 96 vagy 384 mintahelyet (vályút) tartalmazó speciális lemez is elhelyezhető.



2. kép Első nemzedékbeli kiszolgáló-alkatrészek:
1. VA Linux VAR900 2xXeon-500 1 TB felkínált tárhellyel,
2. Raidion.u2w RAID-vezérlők, 3. 2x8x36 GB SCSI lemez és
4. VA Linux 2230-asok és 3x10x72 GB SCSI lemez

A MegaBACE egy SCSI-eszköz, az Applied Biosystems (ABI) 3700 és 3730XL szekvenátorok pedig (6. kép) soros felületen keresztül kezelhetők, az adatokat viszont ethernetkapcsolaton keresztül küldik. Nagy mennyiségű adatot gyűjtenek önműködően, a hozzájuk tartozó program viszont egy mutass és kattints (point-and-click) Windows-alkalmazás. Az ABI gépek a hozzájuk adott helyi Oracle adatbázisba mentik az adatokat. Egy Unix alapú vezérlőprogram forradalmasítaná e gépek kihasználását, különösen a nagyobb laboratóriumokban. Már sikerült csökkentenünk a 3700-es karbantartási munkáit azáltal, hogy az eredetileg a szekvenátorral szállított PC-t IBM x330-as gépre cseréltük (6. kép). A Windows alapú szekvenálórendszernek a linuxos hálózatunkba történő beillesztése remek munka volt az smbmount, az rsync, a Perl és az Apache számára. Az operátor minden egyes sorrend-meghatározási kör befejezésekor beindítja a webvezérlésű adattükörzési folya-



3. kép LIMS sémánk bemutatása. A lemeztábla (sárga) négy táblára hivatkozik (zöld), rá pedig 14 tábla hivatkozik (vörös)



4. kép A vonalkódokat hálózatra kötött Zebra nyomtatók készítik (bal oldalt). A LIMS hordozható felületét iPAQ-ek biztosítják (jobb oldalt)

matot, és az új adatokat a hálózati lemezekre másolja. Tükrözés után az állományokat először a nyers szekvenálójel-formátumból átalakítjuk a tényleges nukleinsavbázisok jelévé és a hozzájuk tartozó minőségértékekké (a meghatározás biztonságának mértéke), majd MySQL-adatbázisban (3.23.55max) tároljuk őket. Ezzel a módszerrel eddig kétféle sorrendet rögzítettünk, azaz körülbelül 1 TB nyers nukleotidsorrend-adatot. A MySQL Laboratory Information Management System (LIMS) adatbázis központi szerepet tölt be a nukleotidsorrend megállapításának folyamataiban. Sémájában 115 táblát, 1171 mezőt és 195 idegen kulcsot találunk. Az adatbázis az összes, a laborral kapcsolatos összetevőt, felszerelést, folyamatot és műveletet követi. Különleges alkalmazáslogika és elnevezési szabályok segítségével sikerült áthidalnunk a MySQL hiányosságát, miszerint nem rendelkezik beépített idegenkulcs-kezeléssel. Az idegen kulcsokat FKTYPE_TABLE__FIELD-nek nevezzük, jelezve, hogy egy TABLE_FIELD-re mutatnak a TABLE táblában. Az idegen kulcs nevének elhagyható TYPE részét akkor használjuk, ha több kulcs



5. kép A laborban szinte minden vonalkóddal van ellátva

is mutat ugyanarra TABLE_FIELD mezőre. A labor szakemberei vonalkódolvasóval kiegészített Wi-Fi Compaq iPAQ gépekkel tartják a kapcsolatot a LIMS adatbázissal (4. kép). Az iPAQ-ok a belső, saját mod_perl készlettel bővített Apache webkiszolgálókra csatlakoznak. A különféle objektumok, azaz a megoldások, lemezek és felszerelések vonalkóddal vannak ellátva (5. kép). A vonalkódokat hálózatra kötött Zebra S600/96XIII vonalkódnymtatóval készítjük nagy ragadóképeségű címkékre (4. kép), amelyek -80 °C (-112 °F) hőmérsékletű hűtőnkben is fennmaradnak. A vonalkódkészítő program Perlben íródott, a címkék formázására a ZPL nyomtatónyelvet használja, a nyomtatást pedig lpr-en keresztül osztja meg. A MegaBACE 1000-es óta laboratóriumunkban a szekvenátorok három nemzedéke fordult már elő, és jelenleg már három ABI 3700-es és három ABI 3730XL gépet (6. kép) üzemeltetünk. A legfrissebb, az ABI 3730XL több 384 mintahelyes lapot képes befogadni, és 1152 DNS minta nukleotidsorrendjét határozza meg 24 óra alatt. Minden egyes minta 700–800, nagy biztonsággal azonosított bázist jelent. Egyetlen 3730XL körülbelül 800 ezer bázist olvas le naponta. A Tor2/SARS genom nukleotidsorrendjének a megállapítását az úgynevezett teljes genomra irányuló (whole-genome shotgun, WGS) módszerrel végeztük. Ennél a megközelítésnél a genom véletlenszerűen kiemelt szakaszait szekvenáljuk redundáns módon, majd utólag állítjuk össze a teljes genomsorrendjét. Tekintve, hogy a szóban forgó vírus méretét körülbelül 30 ezer bázisra becsültük, a teljes genom leolvasásához legalább negyven szekvenciameghatározást kellett végrehajtani. Minthogy azonban a leolvasás véletlen régiókból történt, a minimális olvasásszámnál többet kellett végrehajtanunk, hogy elég átfedésünk legyen a teljes összeállításhoz. A redundancia miatt biztosabbak is lehetünk benne, hogy a genom egyes pozícióin valóban az adott bázist tartalmazó nukleotid áll.



6. kép

Szekvenátorok: 1. MegaBACE 1000, 2. a szekvenátor PC-je, 3. szünetmentes áramforrás, UPS, 4. a szekvenátor áramforrása, 5. ABI 3700-es szekvenátorok, 6. ABI 3730XL és 7. x330 fűrt

A bézs besötétedik

Amikor első IBM x330 kiszolgálóin-
kat vásároltuk, amelyek ma már
egy 168 CPU-t tartalmazó fűrt
részei (7. kép), az 1U felület volt a
kereskedelmi off-the-shelf (COTS)
kategória határa, ahonnan kezdve
élni lehetett a COTS árait. Bézs-
színű gépeinket többé már nem
használnak megosztott számítások-
ra. A komoly terhelésnek alávett
termelési rendszereink, azaz az
Apache és a MySQL, az IBM 4U
x440s-eiben kaptak helyet, ezekben
a nyolcútas hiperszállakkal
(hyperthreading) és 8 GB memó-
riával ellátott Xeon-csomópontok-
ban. A gépeken SuSE 8.1 fut – ez
azon kevés terjesztés egyike,
amelyik képes kezelni az IBM
Summit lapkakészletét. A x440-es
NUMA típusú gép, ahol négypro-
cesszoros modulonként 32 MB L4
gyorstár található, így az IBM
Summit feltja nélkül a rendszer-
mag csak két CPU-t lát. A SuSE
2.4.19 rendszermagja
bigmem+Summit támogatással
mind a nyolc processzor és a 8 GB
memória használatát lehetővé tette.
Ezek az x440-esek még a 2.5-ös
rendszermaghoz viszonyítottan megjel-
elő fejtett NUMA ütemező nélkül is
igen hasznos igavonónak bizonyul-
tak, és lehetővé tették, hogy nyolc
BLAST folyamatot futtassunk pár-
huzamosan, miközben elegendő
memóriánk marad a teljes emberi
genom gyorstárazására a megosz-
tott memóriában. Bárki, aki azt
állítja, hogy a Linux még nem ké-
szült fel a Nagy Vasakra, me-
glepetésre számítson.

Mivel gyorsan növekedtünk, az
NFS alrendszer kezdett problémássá válni. Egészen pontosan
néhány gép összeomlott egy bizonyos NFS kiszolgáló-üggyfél-
váltózat használata esetén. Bár tapasztalataink szerint az

GSC MySQL LIMS

2,1 millió minőségi bázispárt tartalmazó 3250 szekvenciát gyűjtöttünk be, amelyeket a kezdeti
vázlat összeszerkesztéséhez továbbítottunk. Ez körülbelül 70× fedti le redundáns módon a genomot.
A WGS során általában csak 10×-es ismétléssel dolgozunk, de számunkra az időtényező volt a
legfontosabb, így el akartuk kerülni az első sorrend-meghatározási körben nem teljesen lefedett
részek miatt bekövetkező késlekedést.

```
SELECT
SUM(Sequence_Length) AS bp_tot,
AVG(Quality_Length) AS bpq_avg,
SUM(Quality_Length) AS bp_qual_tot,
COUNT(Well) AS reads,
Sequence_DateTime AS date,
Equipment_Name AS equip
FROM
Equipment, Clone_Sequence, Sequence_Batch, Sequence,
Plate, Library, Project
WHERE FK_Sequence_Batch_ID=Sequence_Batch_ID AND
FK_Plate_ID=Plate_ID AND
FK_Library_Name=Library_Name AND
FK_Equipment_ID=Equipment_ID AND
FK_Project_ID=Project_ID AND
FK_Sequence_ID=Sequence_ID AND
Sequence_Subdirectory like "SARS2%" AND
Quality_Length > 100 AND
Sequence_DateTime < "20030413"
GROUP BY Sequence_ID ORDER BY Sequence_DateTime;
```

bp_tot	bpq_avg	bp_tot	reads	date	equip
437256	612.6399	205847	336	2003-04-11 21:07:06	SARS212.B21 D3730-3
412366	752.1074	245187	326	2003-04-11 22:15:34	SARS213.B21 D3730-1
269456	639.1926	225635	353	2003-04-11 22:22:34	SARS215.B21 D3700-6
130525	715.5060	118774	166	2003-04-11 22:25:44	SARS216.B21 D3700-5
282490	682.6311	249843	366	2003-04-11 22:27:14	SARS215.BR D3700-4
310119	612.7601	212015	346	2003-04-11 22:31:56	SARS213.BR D3700-1
182573	681.4975	136981	201	2003-04-11 22:36:40	SARS216.BR D3700-3
301471	642.2273	226064	352	2003-04-12 01:58:16	SARS212.BR D3700-2
401595	690.5204	220276	319	2003-04-12 05:13:26	SARS211.BR D3730-3
460100	642.0468	219580	342	2003-04-12 06:20:52	SARS214.BR D3730-2
182360	471.7832	67465	143	2003-04-12 07:14:44	SARS214.B21 D3730-1

NFS-ügyfelek igen erőteljesek, a jelenlegi Linux NFS szolgálta-
tásokon azért van még mit csiszolni. A leggyorsabb NFS kiszul-
gálónk, egy IBM x342 (2xP3-1.24, 2GB RAM) sem volt képes

A SARS adatainak összeállítására és vizsgálatára az x330-asokat és az x440-est használtuk. A genom nem túl nagy, így az összeállítás egyetlen CPU-n nem vett többet igénybe 15 percnél. Összehasonlításképpen, az emberi genom első nyilvánosságra hozott sorrendje 300 000-szer volt nagyobb a Tor2/SARS méreténél, és az összeállítása négy napon keresztül folyt az UCSC-nél, egy százprocesszoros Linux-fürtön. 2003. április 12-én szombat éjjel 2:25-órakor befejeztük a Tor2/SARS összeállításának hetedik ismétlését, és ezt az állapotot fogadtuk el az első érvényes vázlatként. Ezt importáltuk az AceDB-be, hogy lássuk, mennyire illeszkedik a már ismert proteinkészletekhez (9. kép). A szombatot az összeállításunk kiértékelésével töltöttük, amit aztán egy nappal később az x440-esünk saját, Zope/Plone alapú CMS rendszert futtató nyilvános webkiszolgálójára tettünk fel.

Összegzés

A Tor2/SARS genomját egy negyedik, újfajta coronavírus-csoport tagjaként azonosítottuk, ami információt szolgáltat diagnózisestek, a jövőben pedig esetleges terápia kifejlesztéséhez, beleértve oltóanyag előállítását is. A Linux lehetővé tette, hogy célunkat úgy érjük el, hogy közben nem kell egy vagyont költenünk eszközökre és programokra. Tömegcikként gyártott alkatrészek beépítésével elkerülhettük a hosszú megvalósulási idő miatt bekövetkező értékcsökkenést. Figyelni fogjuk az újonnan felmerülő hibákat, mindeközben MySQL adatbázisunk tárt kapukkal várja az új szekvenciákat.

Köszönetnyilvánítás

A szerzők szeretnének köszönetet mondani *Marco Marra, Steven Jones, Caroline Astell, Rob Holt, Angela Brooks-Wilson,*

Jas Khattra, Jennifer Asano, Sarah Barber, Susanna Chan, Allison Cloutier, Sean Coughlin, Doug Freeman, Noreen Girm, Obi Griffith, Steve Leach, Mike Mayo, Helen McDonald, Steven Montgomery, Pawan Pandoh, Anca Petrescu, Gord Robertson, Jacquie Schein, Asim Siddiqui, Duane Smailus, Jeff Stott és George Yang hölgyeknek és uraknak tudományos szaktudásukért, valamint laboratóriumi és bioinformatikai erőfeszítéseikért. Szeretnénk köszönetet mondani *Kirk Schoeffelnek, Mark Mayonak és Bernard Linek* rendszerfelügyeleti tanácsaiért.

A cikkhez tartozó Kapcsolódó címek az 54. CD Magazin/SARS könyvtárában találhatóak.

Linux Journal 2003. november, 115. szám



Martin Krzywinski (martink@bcgsc.ca)
Bioinformatikai kutató a kanadai Michael Smith Genomtudományi Központban. Idejét fizikai hozzárendelés és adatfeldolgozás-automatizálási kérdések megoldásával tölti a Perl nyelv segítségével.



Yaron Butterfield (ybutterf@bcgsc.ca)
A szekvenáló bioinformatikai csoportot vezeti a kanadai Michael Smith Genomtudományi Központban.

